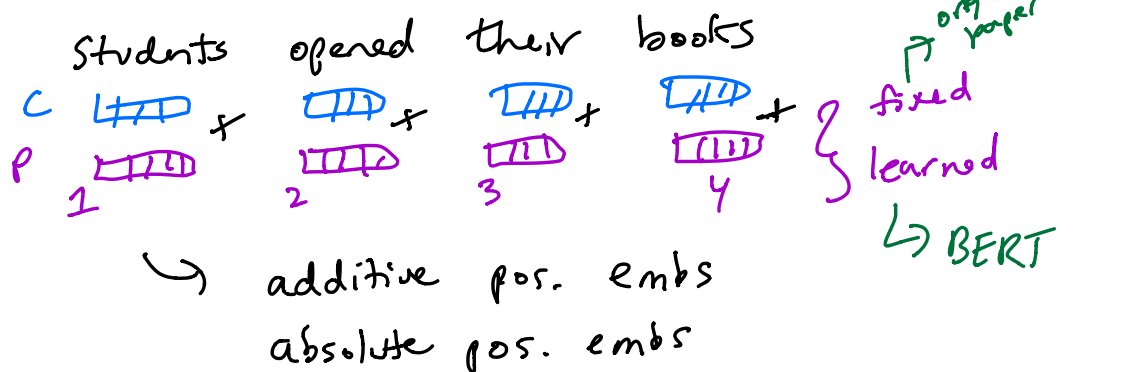


## Position embeddings in Transformers:

- without some explicit injection of position info, self-attention doesn't have any notion of order



absolute vs. relative pos embs

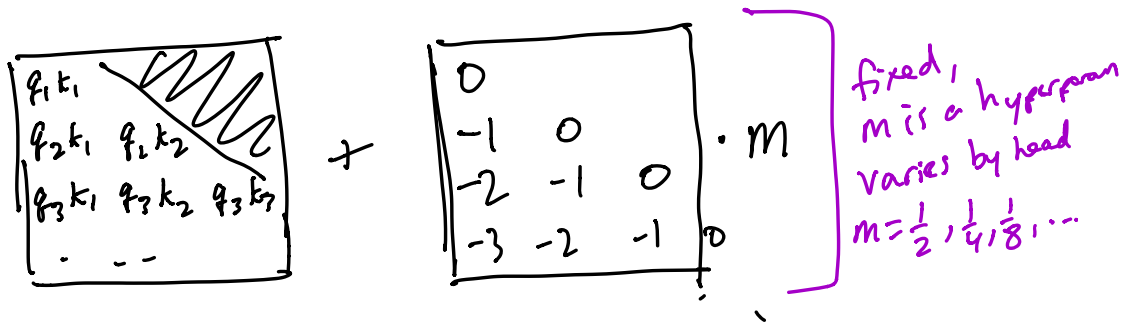
↳ represent every pair of tokens in the input

$$q_{\text{students}} = \underline{W_q \cdot (c_{\text{students}} + p_1)}$$

relative position embs:

- generally cannot be added directly to input embs (RoPE is an exception)
- instead directly modify the attn matrix

ALiBi:  $q_{\text{students}} = f(W_q \cdot C_{\text{students}})$   
 $k_{\text{books}} = f(W_k \cdot C_{\text{books}})$



$\Rightarrow$  intuitively: words that are closer together have a higher dot product

$\Rightarrow$  ALiBi: enables extrapolation beyond the training seq length

$\Rightarrow$  position info is only affecting  $q, k$ , but not  $v$

## Rotary position embs (RoPE)

- enables relative pos. embs without modifying the attn matrix like ALiBi
- instead of adding pos. emb, we actually rotate the  $q, k$  vectors via matrix/vector product w/ a rotation matrix

— goal: dot product of rotated  $q, k$   
 $(q^T k)$  should be a function of  
 relative position only, not abs. pos

ex:  $c_1 \quad c_2 \quad c_3 \quad c_4$   
 students opened their books

we want to compute  $\underline{q_4 \cdot k_1}$

ROPE: find  $f_q, f_k, g$  such that

$$f_q(c_{\text{books}}, 4) = \boxed{\text{||||}} q_4$$

$$f_k(c_{\text{students}}, 1) = \boxed{\text{||||}} k_1$$

$$q_4 \cdot k_1 = g(c_{\text{books}}, c_{\text{students}}, 3) \quad \hookrightarrow 4-1$$

$\Rightarrow$  this can be accomplished by  
 rotating  $W_q c_i$  and  $W_k c_i$  by diff. angles

$$f_q(c_{\text{books}}, 4) = R_{\theta, 4} \cdot W_q c_{\text{books}}$$

$$f_k(c_{\text{students}}, 1) = R_{\theta, 1} \cdot W_k c_{\text{students}}$$

$$g = q^T k$$

$$R_{\theta, t} = \begin{bmatrix} \cos t\theta & -\sin t\theta \\ \sin t\theta & \cos t\theta \end{bmatrix}$$

where  $\theta$  is hyperparameter